

Probabilistic Principal Component Analysis Based on JoyStick Probability Selector

Marko V. Jankovic, *Senior Member, IEEE*, and Masashi Sugiyama

Abstract— Principal component analysis (PCA) is a commonly applied technique for data analysis and processing, e.g. compression or clustering. In this paper we propose a probabilistic PCA model based on the Born rule. In off-line realization it can be seen as a successive optimization problem. In the on-line realization it will be solved by introduction of two different time scales. It will be shown that recently proposed time oriented hierarchical method, used for realization of biologically plausible PCA neural networks, represents a special case of the proposed model. The proposed model gives a general framework for creating different PCA realizations/algorithms. A particular realization can optimize locality of calculation, convergence speed, preciseness or some other parameter of interest. We will present some experimental results to illustrate effectiveness of the proposed model.

I. INTRODUCTION

Principal component analysis (PCA) (e.g. [18]) is a well established technique for dimensionality reduction, and chapters on the subject may be found in numerous texts on multivariate data analysis. Examples of its many applications include data compression, image processing, visualization, exploratory data analysis, pattern recognition, and time series prediction.

The most common derivation of PCA is in terms of a standardized linear projection which maximizes the variance in the projected space [11]. For a set of observed N -dimensional data vectors $\{\mathbf{x}_m\}$, $m \in \{1, \dots, M\}$ the K principal axes \mathbf{w}_k , $k \in \{1, \dots, K\}$, are those orthonormal axes onto which retained variance under projection is maximal. It can be shown that the vectors \mathbf{w}_k are defined by the K dominant eigenvectors \mathbf{u}_k (i.e. those with the largest associated eigenvalues λ_k) of the sample covariance matrix \mathbf{C} , such that $\mathbf{C}\mathbf{u}_k = \lambda_k\mathbf{u}_k$. The K principal components of the observed vector \mathbf{x}_k are given by the vector $\mathbf{y}_k = \mathbf{U}^T(\mathbf{x}_k - \mathbf{x}_{\text{mean}})$, where $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_K)$ and \mathbf{x}_{mean} is the data sample mean. Matrix \mathbf{W} can be expressed as $\mathbf{W} = \mathbf{U}\mathbf{R}$, where \mathbf{R} is a rotation matrix of appropriate dimension. If the variables \mathbf{y}_k are uncorrelated such that the output covariance matrix is diagonal with elements λ_k , then matrix \mathbf{W} can be expressed as $\mathbf{W} = \mathbf{U}\mathbf{P}$, where \mathbf{P} is a permutation matrix.

Artificial neurons and neural networks have been shown to perform PCA when gradient ascent (descent) learning rules are used, which is related to the constrained maximization (minimization) of some objective functions [2, 6, 14, 22]. Due to their low complexity, such algorithms and their implementation in neural networks are potentially useful for tracking of slow changes of correlations in the input data or in updating eigenvectors with new samples.

In recent years, many Principal Subspace Analysis (PSA) algorithms have been proposed and studied in the literature [2, 5, 9, 12, 14, 22]. Some of these PSA algorithms were modified in order to derive parallel PCA algorithms. Usually, the modification was performed by introduction of some asymmetry (inhomogeneity) or nonlinearity in the original PSA algorithm that is not considered desirable from the point of view of the implementation of those algorithms in parallel hardware. Rare examples of a fully homogeneous algorithm are the bigradient algorithm proposed in [19], and the time-oriented hierarchical method proposed in [12, 15]. For a comprehensive review of known parallel, as well as sequential PCA algorithms see e.g. [6].

The Time-Oriented Hierarchical Method (TOHM) is a recently proposed cooperative competitive concept that transforms the PSA algorithms in to PCA algorithms (it can also be used to transform minor subspace algorithms (MSA) into minor component algorithms (MCA)). By implementation of the TOHM, we can create a fully homogeneous learning rule in terms of individual neurons if the initial PSA algorithm is fully homogeneous. The method uses two distinct time scales. A given PSA algorithm is responsible, on a faster time scale, for the behavior of all output neurons. At this scale, a principal subspace is obtained. On a slower time scale, output neurons compete to fulfill their “own interests”. At this scale, basis vectors in the principal subspace are rotated toward the principal eigenvectors. To the best of our knowledge, the bigradient algorithm is the only other known algorithm that uses a similar concept. The difference between the TOHM and the bigradient algorithm has been investigated in [13, 15].

In this paper we will give a probabilistic view of the PCA problem and show that the TOHM is just one possible way to implement PCA on-line. The probabilistic view is based on a, new geometric interpretation of probability based on the Born rule [27]. We will explain that PCA can be seen as a successive optimization problem if we are implementing it off-line. In the case of on-line implementation, we will use constraint optimization on two time scales, for instance as done in the TOHM method. From the aspect of artificial neural networks, the choice of different realization concepts

M. V. Jankovic is with the Department of Computer Science, Tokyo Institute of Technology, Ookayama 2-12-1, Meguro-ku, Tokyo, 152-8552, Japan, on leave from the EE Institute “Nikola Tesla”, Belgrade, Serbia (phone: +81-3-7734-2699; fax: +81-3-7734-2699; e-mail: marko@sg.cs.titech.ac.jp).

M. Sugiyama, is with the Department of Computer Science, Tokyo Institute of Technology, Ookayama 2-12-1, Meguro-ku, Tokyo, 152-8552, Japan, (e-mail: sugi@cs.titech.ac.jp).

has direct impact on the algorithm's convergence speed, preciseness, complexity of plausible hardware realization or biological plausibility. In Section II, the Born rule will be introduced. Geometric interpretation of the Born rule will be given in Section III. Section IV contains introduction of new probabilistic PCA. Some practical aspects of implementation will be discussed in Section V. Experimental results are presented in Section VI. Section VII gives conclusion remarks.

II. QUANTUM PROBABILITY MODEL AND BORN RULE

In quantum mechanics the transition from a deterministic description to a probabilistic one is done using a simple rule termed the Born rule. This rule states that the probability of an outcome (\mathbf{a}) given a state (Ψ) is the square of their inner product $(\mathbf{a}^T \Psi)^2$. This section is based on a similar section in [27].

In quantum mechanics the Born rule is usually taken as one of the axioms. However, this rule has well established foundations. Gleason's theorem [10] states that the Born rule is the only consistent probability distribution for a Hilbert space structure. Wootters [28] has shown that by using the Born rule as a probability rule, the natural Euclidean metrics on a Hilbert space coincides with a natural notion of a statistical distance. Short review for some other justifications of the Born rule can be seen in [27].

The quantum probability model takes place in a Hilbert space H of finite or infinite dimension. A state is represented by a positive semidefinite linear mapping (a matrix ρ) from this space to itself, with a trace of 1, i.e. $\forall \Psi \in H \Psi^T \rho \Psi \geq 0$, $\text{Tr}(\rho) = 1$. Such ρ is self adjoint and is called a density matrix.

Since ρ is self adjoint, its eigenvectors Φ_i are orthonormal and since it is positive semidefinite its eigenvalues p_i are real and nonnegative $p_i \geq 0$. The trace of a matrix is equal to the sum of its eigenvalues, therefore $\sum_i p_i = 1$.

The equality $\rho = \sum_i p_i \Phi_i \Phi_i^T$ is interpreted as "the system is in state Φ_i with probability p_i ". The state ρ is called the pure state if $\exists i$ s.t. $p_i = 1$. In this case, $\rho = \Psi \Psi^T$ for some normalized state vector Ψ , and the system is said to be in state Ψ .

A measurement M with an outcome x in some set X is represented by a collection of positive definite matrices $\{\mathbf{m}_x\}_{x \in X}$ such that $\sum_{x \in X} \mathbf{m}_x = \mathbf{1}$ ($\mathbf{1}$ is being the identity matrix in H). Applying measurement M to state ρ produces outcome x with probability

$$p_x(\rho) = \text{trace}(\rho \mathbf{m}_x) \quad (1)$$

This is the Born rule. Most quantum models deal with a more restrictive type of measurement called the von Neumann measurement, which involves a set of projection operators $\mathbf{m}_a = \mathbf{a} \mathbf{a}^T$ for which $\mathbf{a}^T \mathbf{a}' = \delta_{aa'}$. In a modern language, von Neumann's measurement is a conditional expectation onto a maximal Abelian subalgebra of the algebra of all bounded operators acting on the given Hilbert space. As before, $\sum_{a \in M} \mathbf{a} \mathbf{a}^T = \mathbf{1}$. For this type of measurement the Born rule takes a simpler form: $p_a(\rho) = \mathbf{a}^T \rho \mathbf{a}$. Assuming ρ is a pure state this can be simplified further to

$$p_a(\rho) = (\mathbf{a}^T \Psi)^2. \quad (2)$$

So, we can see that the probability of the outcome of the measurement will be \mathbf{a} , if the state is ρ , is actually cosine square of the angle between vectors \mathbf{a} and Ψ , or $p_a(\rho) = \cos^2(\mathbf{a}, \Psi)$.

III. JOYSTICK PROBABILITY SELECTOR

In this section we will give one novel but simple interpretation of the probability that is related to the Born rule. Here we will assume that we are dealing with finite dimensional discrete variable.

For the moment, let's assume that we are dealing with discrete two dimensional variable. It can associate us with a coin tossing. Assume further that two possible outcomes of our experiment are represented by the dummy variables $\{01\}$ and $\{10\}$. If we represent our coin as a unit norm vector in the two dimensional space (we will call that vector JoyStick Probability Selector or JSPS), then we can have the following simple geometric interpretation given in Fig. 1.

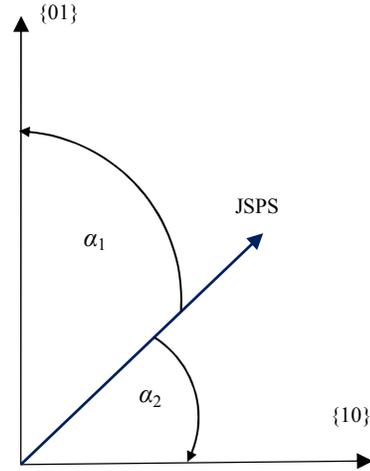


Fig. 1. JoyStick Probability Selector

Now, we will suggest that the probability of outcome $\{01\}$ is equal to cosine square of angle α_1 , while the probability of outcome $\{10\}$ is equal to cosine square of angle α_2 . It is not difficult to see that $\cos^2(\alpha_1) + \cos^2(\alpha_2) = 1$. We can see that the probability of the particular outcome of the experiment (in this case coin toss) is equal to the inner product of the unit norm JSPS and the unit norm vector which represents that outcome. Then, it is easy to see that this coincides to the Born rule interpretation for the case of a pure state and von Neumann measurement system. It is easy to see that when state (JSPS) vector collapses to one of the states, it is not possible to give information of the other state, which is consistent with some quantum mechanics results.

We can check what will happen if our discrete variable is of the dimension 3. In that case our system can be represented by Fig. 2. Now, we have three possible outcomes of the experiment that are represented by dummy variables $\{001\}$, $\{010\}$ and $\{100\}$. Again we have the JSPS vector which

represents the status of our variable before we perform the measurement. Again, the probability of the outcome is given by the cosine square of the angle between JSPS and particular

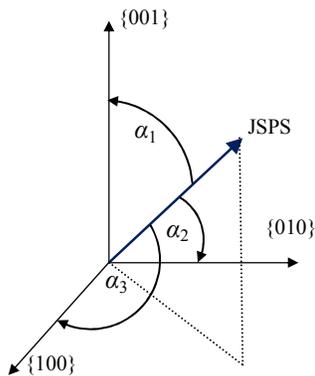


Fig. 2. A 3-D example

outcome vector. It is not difficult to check that

$$\sum_{i=1}^3 \cos^2(\alpha_i) = 1. \quad (3)$$

It can be easily concluded that it follows from Pythagorean Theorem. For any 3-D vector whose norm is r we have

$$\sum_{i=1}^3 r^2 \cos^2(\alpha_i) = r^2. \quad (4)$$

It can be easily concluded that this way of reasoning can be extended to any finite dimension D .

It is interesting to notice that this representation is actually reasonable for unimodal data (data that can be represented by a pure state). If the data is multi modal we need more complex model which is not going to be analyzed here.

IV. PROBABILISTIC PCA

PCA also has another convenient interpretation that is perhaps less well known. In this probabilistic interpretation [24], each point \mathbf{x}_i is thought of as a random draw from some unknown distribution P_{θ_i} , where P_{θ} denotes a unit Gaussian with mean $\theta \in \mathbb{R}^N$. Then, the purpose of PCA is to find the set of parameters $\theta_1, \dots, \theta_m$ that maximize the likelihood of the data, subject to the condition that these parameters all lie in a low-dimensional subspace. In other words, $\mathbf{x}_1, \dots, \mathbf{x}_m$ are considered to be a noise-corrupted version of some true points $\theta_1, \dots, \theta_m$ which lie in the subspace; the goal is to find these true points, and the main assumption is that the noise is Gaussian. The equivalence of this interpretation to the ones based on variance maximization follows simply from the fact that negative log likelihood under this Gaussian model is equal (ignoring constants) to the sum of the squared distances from the data points \mathbf{x}_i to their projections θ_i in the subspace. The Gaussian assumption may be inappropriate, for instance, if the data is binary-valued, or integer-valued, or is nonnegative [7]. In fact, the Gaussian is only one of the canonical distributions that make up the exponential family, and it is a distribution tailored to real-valued data. The

Poisson is better suited to integer data and Bernoulli to binary data. It seems natural to consider variants of PCA which are founded upon these other distributions in place of the Gaussian.

So, in [7] PCA was extended to the rest of the exponential family. Let $\{P_{\theta}\}$ be any parameterized set of distributions from the exponential family, where θ is the natural parameter of a distribution. Given data $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^N$, the goal is now to find parameters $\theta_1, \dots, \theta_m$ which lie in a low dimensional subspace and for which the log-likelihood is maximized.

That approach effortlessly permits hybrid dimensionality reduction schemes in which different types of distributions can be used for different attributes of the data. If the data \mathbf{x}_i have few binary attributes and a few integer-valued attributes, then some coordinates of the corresponding θ_i can be the parameters of binomial distributions while others are the parameters of Poisson distributions.

In [7] was pointed out that dimensionality reduction schemes for non-Gaussian distributions are substantially different from PCA. For instance, in PCA the parameters θ_i , which are means of Gaussians, lie in a space which coincides with that of the data \mathbf{x}_i . This is not the case in general, and therefore, although the parameters θ_i lie in a linear subspace, they typically correspond to a nonlinear surface in the space of the data.

Here, we have to point out two things when we are considering the existing probabilistic PCA approaches:

- First, it is frequently named probabilistic PCA although it should be correctly named probabilistic PSA. Although, PSA and PCA are similar techniques they should be considered as different. First of all, result of PCA takes a solution from a discrete set while the result of PSA is taken from a continuous set. Another important difference is that the input data covariance matrix is diagonal only in the coordinate system which is defined by eigenvectors themselves and not in any other rotated system.
- The second problem is that proposed probabilistic models are more related to factor analysis [7] and make strong assumptions about the distribution of the noise such as Gaussian or in more general case exponential family. It seems that those assumptions are too strong and as we are going to see, not really necessary.

Now, we will give a new definition of PCA that is different from already mentioned definitions that actually specify PSA. The definition is based on modification of the Hotteling definition [11]:

Definition 1: PCA can be defined as a problem of two successive optimization processes. In the first step, PSA is performed – which means that we are going to find orthogonal projection of the data onto lower dimensional linear space, known as principal subspace, such that the projected data is maximized. (In the text that follows, this step will have a different interpretation based on the Born rule, but we will avoid it here in order to simplify the definition.) In the second step the data projected in the principal subspace will be optimized in the sense that the entropy of the data is

minimized. In the second step, we perform learning on Grassman Principal Submanifold [13].

This formulation can be related to the recently proposed Time Oriented Hierarchical Method (TOHM) [12, 15] which represents a cooperative-competitive concept and is known as a general method that transforms PSA/MSA algorithms into PCA/MCA algorithms. The TOHM enables on-line realization in a single step (which means the data can be used only once; the problem is solved by introduction of multiple time scales) and will be presented later in the paper.

Formulation can also be seen as an approach similar to the one proposed in [8]. In this case first we are looking for a pool of the possible solutions (hyperensemble) and after we define it we are looking for the solution (ensemble) that fulfils the additional criterion. In other words, we are breaking one difficult problem into two simpler problems.

It is not difficult to see that the probabilistic MCA can be defined analogously as a two step successive optimization process. In the first step, the data is projected in the minor subspace, and in the second step entropy of the projected data should be minimized.

In the text that follows, we will concentrate on symmetric cost functions in the sense that possible artificial neural network realization is fully symmetric from the point of view of an individual neuron. Of course, it is not very difficult to give definitions that are more general so that they can include nonsymmetric cost functions (e.g. to cover the cases proposed in [1] or [16]), but it will not be discussed in this paper. In this paper, we consider nonsymmetric realizations as not completely successful realizations of the fully symmetric case, which is opposite to the point of view in which the fully symmetric realization represents a special case of nonsymmetric realizations.

Based on the Definition 1 and quantum probabilistic model we are going to give definition of probabilistic model for PCA. In order to do that, first we are going to give a new probabilistic approach to PSA.

A. Probabilistic PSA

Without a loss of generality, let's assume that we are dealing with a random variable \mathbf{x} that take realizations from a set of observed N -dimensional zero-mean data vectors $\{\mathbf{x}_i\}$, $i \in \{1, \dots, N_{\text{sample}}\}$ which are sampled from some distribution in time instants $t = iT$ where i is already defined and T represents the sampling period. We will denote $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$ the column vectors of orthonormal matrix \mathbf{W} ($K \leq N$). Now, we are going to define probabilistic PSA. In order to do that, we are going to specify the following cost function

$$\text{JS1}(\mathbf{W}) = \text{E}(p(\mathbf{W}, \mathbf{x})) = \sum_i p(\mathbf{W}, \mathbf{x}_i), \quad (5)$$

where E represents the empirical expectation operator, \mathbf{W} is orthogonal matrix and $p(\mathbf{W}, \mathbf{x})$ is defined in the text that follows. This formulation represents only a class of algorithms that explicitly require orthonormality of matrix \mathbf{W} . For the cases where explicit orthonormality is not desirable constraint (e.g. in biologically plausible neural networks – since the explicit orthonormality usually requires global computation) we can formulate a different cost function

which requires explicitly only unit norm of columns of the matrix \mathbf{W} (that can be obtained by local calculations). This is the case for MH(O) type learning algorithms [14]. Here it will not be done due to space limitation.

We can write

$$p(\mathbf{W}, \mathbf{x}_i) = p(\mathbf{W}|\mathbf{x}_i)p(\mathbf{x}_i), \quad (6)$$

where we define

$$p(\mathbf{W}|\mathbf{x}_i) = \sum_{k=1}^K p(\mathbf{w}_k|\mathbf{x}_i). \quad (7)$$

Now, we will define $p(\mathbf{w}_k|\mathbf{x}_i)$ by usage of the Born rule. In that case we have

$$p(\mathbf{w}_k|\mathbf{x}_i) = \frac{(\mathbf{w}_k^T \mathbf{x}_i)^2}{\|\mathbf{x}_i\|^2}. \quad (8)$$

In other words $p(\mathbf{w}_k|\mathbf{x}_i) = \cos^2(\mathbf{w}_k, \mathbf{x}_i)$. Generally speaking, we consider the data as a state vectors in the von Neuman measurement system which is defined by the columns of the matrix \mathbf{W} .

Also, we are going to define $p(\mathbf{x} = \mathbf{x}_i | t=iT)$ as

$$p(\mathbf{x} = \mathbf{x}_i) = \frac{\|\mathbf{x}_i\|^2}{\sum_{n=1}^{N_{\text{sample}}} \|\mathbf{x}_n\|^2}, \quad (9)$$

where N_{sample} represents the overall number of samples that are going to be analyzed. It is interesting to note that the only thing that we can conclude about the $p(\mathbf{x}_i)$ is that it is proportional to $\|\mathbf{x}_i\|^2$. The sum in the denominator represents the energy of samples that are going to be analyzed – we actually do not know the value of that sum in any, but the final moment. However, we know that it represents some constant. We can easily see that the adopted probability measure fulfils the two conditions that are required for probability function $f(\mathbf{x})$ (in our case $p(\mathbf{x})$) to be considered as a modified generalized probability measure [26]:

1. For each \mathbf{x} , $0 \leq f(\mathbf{x}) \leq 1$,
2. $\sum_i f(\mathbf{x}_i) = 1$.

If we want to calculate marginal probability $p(\mathbf{x} = \mathbf{x}_s)$ (i.e. probability of observing a certain value \mathbf{x}_s no matter what is the time instant i) we can do it as

$$p(\mathbf{x} = \mathbf{x}_s) = \sum_{i=1}^{N_{\text{sample}}} \delta(\mathbf{x}_i, \mathbf{x}_s) p(\mathbf{x}_i), \quad (10)$$

where δ represents Kronecker δ function. In some sense, our approach is similar to the Dirichlet process model which is usually used to model a stream of symbols where the vocabulary of symbols is not limited [20]. In the Dirichlet process case we have one parameter α which is constant. In our case, it is not constant but it is defined by the energy of the new incoming symbol (sample). Also, in our case every incoming symbol is treated as a new one – so initially we do not consider that the same samples in different moments are the same (e.g. we can make analogy to spin of particle). That

means that our input space is initially considered as an augmented space of dimension N_{sample} . When we observe all samples, we can switch to the “deformed” space of dimension N_d , which is defined by the number of different energy levels d .

We can see that (5) can be written as

$$\begin{aligned} JS1(\mathbf{W}) &= E(\mathbf{W}, \mathbf{x}) = \sum_i p(\mathbf{W}, \mathbf{x}_i) \\ &= \frac{1}{\sum_{j=1}^{N_{sample}} \|\mathbf{x}_j\|^2} \left(\sum_{k=1}^K \sum_{i=1}^{N_{sample}} \mathbf{w}_k^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{w}_k \right). \end{aligned} \quad (11)$$

Now, we define probabilistic PSA as a constrained maximization problem

$$\max_{\mathbf{W}} JS1(\mathbf{W}), \quad (12)$$

under the constraint that \mathbf{W} is an orthonormal matrix. So, we are looking for \mathbf{W} which maximizes joint probability $p(\mathbf{W}, \mathbf{x})$ on the average. Solution will be obtained for matrix $\mathbf{W}=\mathbf{W}^{JS1}$ which represents the rotated matrix of matrix \mathbf{U} which consists of eigenvectors \mathbf{u}_i as columns, where \mathbf{u}_i represent K dominant eigenvectors of the covariance matrix $\mathbf{C} = E(\mathbf{x}\mathbf{x}^T)$. Again, the matrix \mathbf{W} is selected in such a way that we maximize, on the average, joint probability $p(\mathbf{W}, \mathbf{x})$, where the columns of \mathbf{W} represent measurement system with \mathbf{x} as a state vectors. It is not difficult to see that with this definition, PSA also can be extended to a clustering problem – see e.g. [17, 27]. Also, it is easy to see that application of gradient ascent to JS1 maximization will lead to Hebb learning rule.

B. Probabilistic PCA

So, as a result of PSA we will obtain a “pool” or “hyperensemble” of matrices \mathbf{W}^{JS1} that span the subspace of the empirical covariance matrix of zero-mean input samples \mathbf{x} . Now, in order to create PCA, we are going to do the second step. In this step we are going to select an ensemble of solutions that fulfill additional criterion, once the matrix \mathbf{W} satisfies the solution of PSA problem.

Let’s assume that the matrix $\mathbf{W}=\mathbf{W}^{JS1}$ is an orthonormal matrix and that it spans the principal subspace. Then, we are going to analyze the following cost function

$$JS2(\mathbf{W}^{JS1}) = E\left(S(p(\mathbf{W}^{JS1}, \mathbf{x}))\right) \quad (13)$$

where S represents the entropy (see next paragraph). If we chose the Tsallis entropy [25], (13) can be written as

$$JS2(\mathbf{W}^{JS1}) = E\left(\frac{1 - \sum_{k=1}^K p^q(\mathbf{w}_k^{JS1}, \mathbf{x})}{q-1}\right), \quad (14)$$

where q is a nonnegative real value, different from 1, and again E represents empirical expectation. In the case $q \rightarrow 1$, we will have the Shannon entropy $S=S_S$, so our cost function is defined as

$$JS2(\mathbf{W}^{JS1}) = E\left(\sum_{k=1}^K p(\mathbf{w}_k^{JS1}, \mathbf{x}) \log\left(\frac{1}{p(\mathbf{w}_k^{JS1}, \mathbf{x})}\right)\right). \quad (15)$$

We can, also, use some other definition of entropy. Now, probabilistic PCA can be defined as the following constrained minimization problem

$$\min_{\mathbf{W}^{JS1}} JS2(\mathbf{W}^{JS1}), \quad (16)$$

under constraint that \mathbf{W}^{JS1} is orthonormal and spans the principal subspace of matrix \mathbf{C} . It can be shown that the solution of this optimization process will be $\mathbf{W}^{JS1} = \mathbf{U}\mathbf{P}$, where, as before, \mathbf{U} consists of K principal eigenvectors \mathbf{u}_i as columns, \mathbf{u}_i represent principal eigenvectors of the covariance matrix $\mathbf{C} = E(\mathbf{x}\mathbf{x}^T)$, and \mathbf{P} is a permutation matrix. This can be seen from section 11.3 in [21] for the Shannon entropy and can be done in a similar way for other types of entropy (proof is omitted here, due to space limitation).

We will have to emphasize that, strictly speaking, proposed equations are correct only in the case $K=N$, since in the case $K < N$, $\sum_{k=1}^K p(\mathbf{w}_k^{JS1}, \mathbf{x}) < 1$. In the case $K < N$ we will define the term “principal subspace” entropy, S^{PS} . Here, we will analyze only the Shannon entropy case, but similar can be done for Tsallis or other types of entropy. We define

$$S^{PS}(\mathbf{W}^{JS1}) = \sum_{k=1}^K p(\mathbf{w}_k^{JS1}, \mathbf{x}) \log\left(\frac{1}{p(\mathbf{w}_k^{JS1}, \mathbf{x})}\right). \quad (17)$$

Correct definition of Shannon entropy in the principal subspace can be done as

$$S_S(\mathbf{W}^{JS1}) = \sum_{k=1}^K p^{PS}(\mathbf{w}_k^{JS1}, \mathbf{x}) \log\left(\frac{1}{p^{PS}(\mathbf{w}_k^{JS1}, \mathbf{x})}\right) \quad (18)$$

where

$$p^{PS}(\mathbf{w}_k^{JS1}, \mathbf{x}) = \frac{p(\mathbf{w}_k^{JS1}, \mathbf{x})}{\sum_{i=1}^K p(\mathbf{w}_i^{JS1}, \mathbf{x})} = \frac{p(\mathbf{w}_k^{JS1}, \mathbf{x})}{Cl}. \quad (19)$$

Definition of constant Cl is obvious from (19), and now we have $\sum_{i=1}^K p^{PS}(\mathbf{w}_i^{JS1}, \mathbf{x}) = 1$. Now, we can write (18) as

$$S_S(\mathbf{W}^{JS1}) = \log(Cl) + \frac{1}{Cl} S^{SP}(\mathbf{W}^{JS1}). \quad (20)$$

From (20) can be seen that minimization of $S_S(\mathbf{W}^{JS1})$ is equivalent to minimization of $S^{SP}(\mathbf{W}^{JS1})$, so equations (14) and (15) can be formally used even in the cases when $K < N$. Other way to show this easily, would be to project input data into principal subspace and then to use such data as an input data in minimization of JS2.

Equation (13) can be made more general if we use a monotonic increasing function of entropy instead of the entropy itself.

C. Minimization of approximate Tsallis entropy

In this section we will give a very simple derivation of the special case of the proposed PCA method, which will be also useful when we, later, discuss how PCA can be realized on-line. In a special case, PCA can be seen as a minimization of the approximate Tsallis entropy, $S_T(q \approx 1)$.

Let's consider the Tsallis entropy with parameter $q = 1 + \varepsilon$, where ε is very small. In that case we have

$$E(S_T(\mathbf{W}, q)) = E\left(\frac{1 - \sum_{k=1}^K p^q(\mathbf{w}_k, \mathbf{x})}{q-1}\right). \quad (21)$$

When q is close to 1 we can have the following approximation

$$p^{q-1}(\mathbf{w}_k, \mathbf{x}) = e^{(q-1)\ln(p(\mathbf{w}_k, \mathbf{x}))} \approx 1 + (q-1)\ln(p(\mathbf{w}_k, \mathbf{x})). \quad (22)$$

Then, we can write

$$E(S_T(\mathbf{W}, q)) \approx E\left(\frac{1 - \sum_{k=1}^K p(\mathbf{w}_k, \mathbf{x})}{q-1} - \sum_{k=1}^K p(\mathbf{w}_k, \mathbf{x}) \ln(p(\mathbf{w}_k, \mathbf{x}))\right), \quad (23)$$

or

$$E(S_T(\mathbf{W}, q)) \approx E\left(\frac{1 - \sum_{k=1}^K p(\mathbf{w}_k, \mathbf{x})}{q-1}\right) + E(S_S(\mathbf{W})). \quad (24)$$

So, in this case, when q is close to 1, the Tsallis entropy can be approximately expressed as the sum of two terms – one that depends on joint probabilities and the other that represents the Shannon's entropy. The first term is multiplied by some factor $(1/(q-1))$ which is much bigger than 1, like it is calculated on a faster time-scale than the second term that represents entropy. Now, based on discussion in previous two subsections, we can conclude that minimization of the approximate Tsallis entropy will lead toward PCA. It is interesting to notice that in the case $N=K$ we only need to minimize the Shannon entropy term since the $\sum_j p(\mathbf{w}_j, \mathbf{x}) = 1$. Also, we point out that in on-line learning algorithms, where the matrix \mathbf{W} is only approximately orthonormal (since the norms of the columns of the matrix \mathbf{W} are slightly different from 1), $\sum_j p(\mathbf{w}_j, \mathbf{x})$ can be understood more precisely as $\sum_j p(\mathbf{w}_j, \mathbf{x})^a$, for some a that is not exactly one (so the effective exponent is actually a little bit different from 1). In the case of strongly orthonormally constrained (SOC) algorithms, which means that \mathbf{W} is orthonormal in all learning steps, effective exponent $a = 1$.

V. PRACTICAL ASPECTS OF IMPLEMENTATION

In this section we will discuss some practical aspects of the implementation of the proposed method. We will discuss off-line implementation and on-line implementation (samples can be used only once) cases separately.

First, based on discussion in the previous section, we can sketch a following algorithm:

Algorithm for off-line calculation of PCA:

In: Matrix \mathbf{X} of size $N \times N_{sample}$

In: Desired number of principle components

Computation:

Step 1:

Compute $\mathbf{W}^{JS1} = \max_{\mathbf{W}} JS1(\mathbf{W})$ where

$$JS1(\mathbf{W}) = E(\mathbf{W}, \mathbf{x}) = \sum_i p(\mathbf{W}, \mathbf{x}_i) \\ = \frac{1}{\sum_{j=1}^{N_{sample}} \|\mathbf{x}_j\|^2} \left(\sum_{k=1}^K \sum_{i=1}^{N_{sample}} \mathbf{w}_k^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{w}_k \right)$$

Step 2:

Compute $\mathbf{W}^{JS2} = \min_{\mathbf{W}^{JS1}} JS2(\mathbf{W}^{JS1})$ where

$$JS2(\mathbf{W}^{JS1}) = E\left(S\left(p(\mathbf{W}^{JS1}, \mathbf{x})\right)\right).$$

From the implementation point of view, this realization is not optimal because we should find the pool of the matrices \mathbf{W}^{JS1} that minimize function JS1. This is impossible if we want to find all of them and can be pretty demanding from the point of view of calculation time and storage capacity if we want to have a big, but finite number of such solutions. Instead, it is much better to find one solution \mathbf{W}^{JS1*} , and to use it to project input data \mathbf{X} to \mathbf{X}^{PSA} that lies in the principal subspace. In the first step we can use any of the available algorithms for PSA calculation. Then we can use the \mathbf{X}^{PSA} as an input to the Step 2 of the algorithm, and find a solution \mathbf{W}^{JS2} that minimizes the cost function $JS2(\mathbf{W}^{JS1*})$. In that case the algorithm looks like:

Improved algorithm for off-line calculation of PCA:

Input:

In: Matrix \mathbf{X} of size $N \times N_{sample}$

In: Desired number of principle components K

Computation:

Step 1:

Compute \mathbf{W}^{JS1} that spans the principal subspace of the empirical covariance matrix $\mathbf{C} = \mathbf{X} * \mathbf{X}^T$. Then, project \mathbf{X} to that subspace and get $\mathbf{X}^{PSA} = \mathbf{W}^{JS1T} * \mathbf{X}$.

Step 2:

Use \mathbf{X}^{PSA} as an input matrix and select some entropy measure S .

Compute $\mathbf{W}^{JS2} = \min_{\mathbf{W}^{JS1}} JS2(\mathbf{W}^{JS1})$ where

$$JS2(\mathbf{W}^{JS1}) = E\left(S\left(p(\mathbf{W}^{JS1}, \mathbf{x}^{PSA})\right)\right).$$

In realization of step 1, we can use any existing off-line PSA learning algorithm, while in step 2 we can use any algorithm that minimizes selected entropy function.

In on-line realization, we can use every sample only once – which means we cannot use the successive optimization

technique that is proposed. The solution can be based on minimization of approximate the Tsallis entropy already explained in subsection IV.C. Here it will be introduced as a cooperative-competitive concept that was used in artificial neural networks and which is called the time oriented hierarchical learning (TOHM). In this approach, two time scales are proposed (it can be easily modified for the case of multiple-time scales). The idea of multiple time scales is used in physics for simulation purposes [4, 23], and, also, represents the major idea that support the superstatistics concept [3].

In TOHM on a faster time-scale we calculate parameters of global interest (principal subspace), while on the slower time scale, we calculate parameters that are of local interest (principal eigenvectors). Concrete realization is based on Lagrangian multiplier idea – solution that should be obtained on the faster time scale is used as a constraint on the slower time-scale. This can be formally represented as the maximization of the following cost function

$$JS(\mathbf{W}) = JS1(\mathbf{W}) - \mu JS2(\mathbf{W}) \quad (25)$$

where μ is a real number which is by absolute value smaller than one (since the JS2 should be realized on a slower time-scale). We can see here that Lagrangian multiplier is used in an unusual way – it does not multiply the constraint part, but the objective one. Anyway, it will not affect the solution since it just represents the weight ratio of the different parts of the learning rule.

Several algorithms that perform PCA, and that are based on the TOHM method are proposed in the literature [12], [13], [15] and [16]. Without going into the details, in all those approaches the PCA algorithms are obtained from some known PSA algorithms together with minimization of the Tsallis entropy, $S_T(q \approx 1)$. Value of μ was selected in all algorithms independently, after stability analysis for every particular choice of JS1 and JS2.

VI. EXPERIMENTAL RESULTS

In this section, we will present some small scale simulation results. We will consider a single layer linear artificial neural network with four inputs and two outputs. In all experiments, the adopted PSA algorithm was as in [22] and PCA parts of algorithms were obtained by minimization of four different entropy functions, or more precisely, the principal subspace entropy functions. Selected entropy functions were Tsallis entropy with $q=0.5$, $q \approx 1$ and $q=2$ as well as the Shannon entropy function. In figures that follow, Error was defined as

$$\text{Error} = \log(\text{abs}(\text{abs}(\cos(\mathbf{v}_1, \mathbf{u}_1)) - 1))$$

where \mathbf{v}_1 represents principal eigenvector calculated by neural network, while \mathbf{u}_1 represents numerically calculated principal eigenvector of input covariance matrix. In all experiments, the learning factors ($\sim 1/t$) and the initial value of the weight matrix were the same.

In Fig. 3, the input signals were generated as a linear mixture of random Gaussian noise. In Fig. 4 input signals

were generated as linear mixture of random uniform noise. Sinusoidal signals were used as input in Fig. 5. Input signals in Fig. 6 were generated as a mixture random uniform noise but learning rate is decreased 10 times.

We can see that, depending on the type of input signals, we have different behaviors of algorithms. Although we cannot say that the choice of the particular entropy measure will

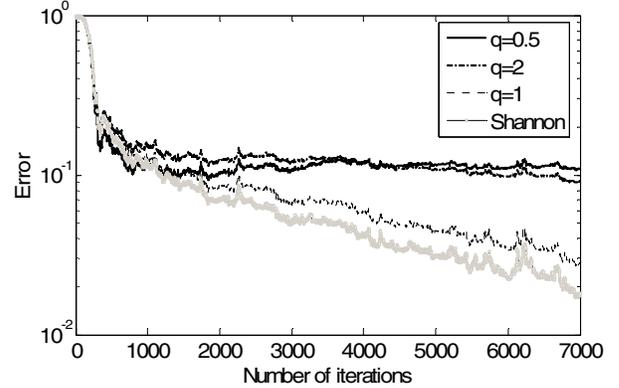


Fig. 3. Input is linear mixture of random Gaussians signals

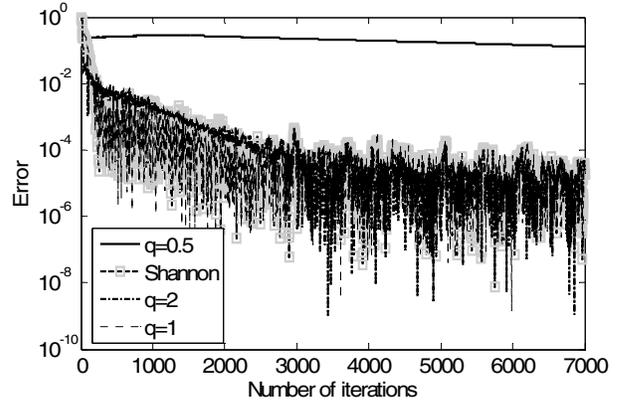


Fig. 4. Input is a linear mixture of random uniformly distributed signals

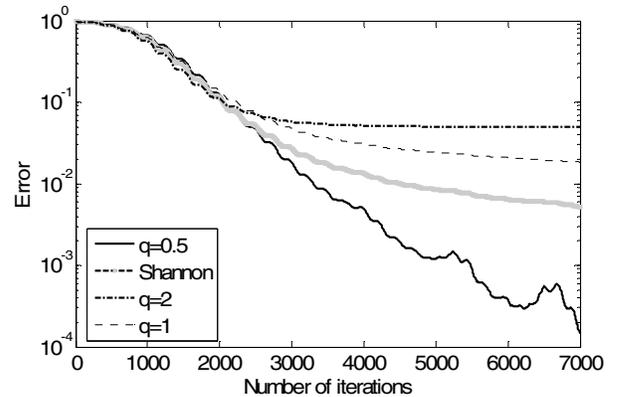


Fig. 5. Input is represented by sinusoidal signals

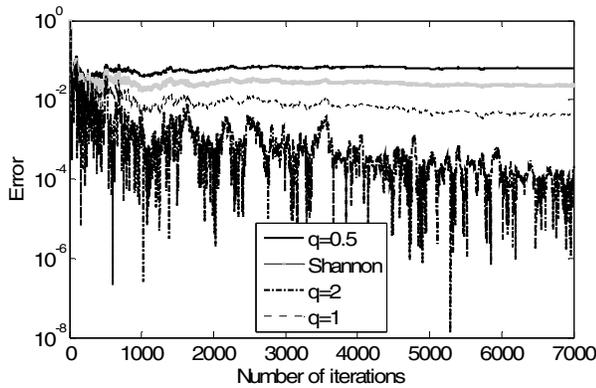


Fig. 6. Input is same as in a picture 4, but learning rate is decreased 10 times

always lead to fastest convergence and highest precision, based on presented experimental results, we can see that selection of the Shannon entropy and the Tsallis entropy with $q \approx 1$, will mainly result in a satisfactory performance. In some sense, those two choices represent robust/general entropy measures, while the two others represent entropy measures suitable for signals or networks with specific features.

VII. CONCLUSION

In this paper we proposed a probabilistic principal component analysis based on the Born rule and its new, simple geometrical representation, named JoyStick Probability Selector. Proposed probabilistic PCA can be solved off-line as a successive optimization problem. In on-line implementation, the problem is solved by the introduction of two distinct time-scales. By selection of different PSA algorithms and different entropy functions, it is possible to create a big number of algorithms that can be used for PCA. In that way, it is possible to create algorithms that could be potentially optimal from the point of view of convergence speed, preciseness, complexity of hardware implementation, locality of calculations etc. A few simulation results were used to illustrate how the choice of entropy function affects the convergence speed and preciseness of the PCA algorithm implemented in an artificial neural network setup.

Proposed probabilistic PSA was based on existence of explicit orthonormality constraint for weight matrix. In the papers that follow, it will be explained how proposed model can be modified to include the cases where the orthonormality constraint is implicit (e.g. MH(O) type algorithms).

As it is briefly mentioned, the proposed approach can be easily modified to give definition of probabilistic MCA.

It is possible to show that similar way of reasoning can lead to a new, probabilistic definition of independent component analysis (ICA). By doing this, it is possible to give different interpretations of some existing ICA solutions and to propose new ones. This is currently under investigation.

ACKNOWLEDGMENT

The authors are grateful to Dr. Ryota Tomioka for useful discussions which improved paper's quality and readability.

REFERENCES

- [1] J.-H. Ahn, J.-H. Oh and S. Choi, "Learning principal directions: Integrated-Squared-Error minimization", *Neurocomputing*, vol. 70, pp. 1372- 1381, 2007.
- [2] P. Baldi and K. Hornik, "Learning in linear neural networks: A survey", *IEEE Trans. Neural Networks*, vol. 6, pp. 837-858, 1995.
- [3] C. Beck, "Superstatistics: theory and applications", *Continuum Mech. Thermodyn.*, vol. 16, pp. 293-304, 2004.
- [4] C.M. Bishop, *Pattern Recognition and Machine Learning*. New York: Springer Science+Business Media, LLC, 2006.
- [5] T.-P. Chen and S.-I Amari, "Unified stabilization approach to principal and minor components", *Neural Networks*, vol. 14, pp. 1377-1387, 2001.
- [6] A. Cichocki and S.-I. Amari, *Adaptive Blind Signal and Image Processing – Learning Algorithms and Applications*. John Wiley and Sons, 2003.
- [7] M. Collins, S. Dasgupta, and R. E. Schapire, "A generalization of principal component analysis to the exponential family", *Advances in Neural Information Processing Systems*, vol. 14, pp. 617-624, 2002.
- [8] G. E. Crooks, "Beyond Boltzmann-Gibbs statistics: Maximum entropy hyperensembles out of equilibrium", *Physical Review*, E 75: 041119, 2007.
- [9] P. Földiák, "Adaptive network for optimal linear feature extraction", *IJCNN'89*, pp. 401-405, 1989.
- [10] A. M. Gleason, "Measure on the closed subspaces of a Hilbert space", *J. Math. & Mechanics*, vol. 6, pp. 885-893, 1957.
- [11] H. Hotelling, "Analysis of a complex of statistical variables into principal components", *Journal of Educational Psychology*, vol. 24, pp. 417-441, 1933.
- [12] M. Jankovic and H. Ogawa, "Time-Oriented hierarchical method for computation of principal components using subspace learning algorithm", *Int. J. Neural Systems*, vol. 14, pp. 313-324, 2004.
- [13] M. Jankovic and B. Reljin "Neural learning on Grassman/Stiefel principal/minor submanifold", EUROCON 2005, pp. 249-252, 2005.
- [14] M. Jankovic and H. Ogawa, "Modulated Hebb-Oja learning rule – A method for principal subspace analysis", *IEEE Trans. on Neural Networks*, vol. 17, pp. 345-356, 2006.
- [15] M. Jankovic and B. Reljin, "General stochastic approximation and neural learning on principal Stiefel submanifold", *WSEAS Trans. on Information Science and Applications*, vol. 3, pp. 1195-1201, 2006.
- [16] M. Jankovic and M. Sugiyama, "A multipurpose linear component analysis method based on modulated Hebb Oja learning rule", *IEEE Signal Processing Letters*, vol. 15, pp. 677-680, 2008.
- [17] M. Jankovic, M. Sugiyama and B. Reljin, "Tensor based image segmentation", in *Proc. 9th Symposium on Neural Network Applications in Electrical Engineering*, Belgrade, 2008, pp. 67-73.
- [18] I. T. Jolliffe, *Principal Component Analysis*, 2nd ed., New York: Springer-Verlag, 2002.
- [19] J. Karhunen, E. Oja, L. Wang, R. Vigarío and J. Joutsalo, "A class of neural networks for independent component analysis", *IEEE Trans. on Neural Networks*, vol. 8, pp. 486-504, 1997.
- [20] D. J. C. MacKay, *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, 2003.
- [21] M. A. Nielsen and I. L. Chuang, *Quantum Computation and Quantum Information*. Cambridge University Press, 2000.
- [22] E. Oja, "Neural Networks, principal components, and subspaces", *Int. J. Neural Syst.*, vol. 1, pp. 61-68, 1989.
- [23] J. C. Sexton and D. H. Weingarten, "Hamiltonian evolution for the Hybrid Monte Carlo algorithm", *Nuclear Physics B*, vol. 380, pp. 665-677, 1992.
- [24] M. Tipping and C. M. Bishop, "Probabilistic Principal Component Analysis", *Journal of the Royal Statistical Society, Series B*, vol.61, part 3, pp. 611-622, 1999.
- [25] C. Tsallis, "Possible generalization of Boltzmann-Gibbs statistics", *J. Statistical Physics*, vol. 52, pp. 479-478, 1988.
- [26] M. K. Warmuth and D. Kuzmin, "A Bayesian probability calculus for density matrices", in *Proc. 22nd Conference on Uncertainty in Artificial Intelligence*, 2006, pp. 503-511.
- [27] L. Wolf, "Learning using the Born rule", Tech. Rep. MIT-CSAIL-TR-2006-036, 2006.
- [28] W. K. Wothers, "Statistical distance and Hilbert space", *Physical Review D.*, vol. 23, pp. 357-362, 1981.